

Une approche de construction d'espaces de représentation multidimensionnels dédiés à la visualisation

Riadh Ben Messaoud, Kamel Aouiche, Cécile Favre

Laboratoire ERIC, Université Lumière Lyon 2

5 avenue Pierre Mendès-France

69676 Bron Cedex

{rbenmessaoud | kaouiche | cfavre}@eric.univ-lyon2.fr

Résumé. Dans un système décisionnel, la composante visuelle est importante pour l'analyse en ligne OLAP. Dans cet article, nous proposons une nouvelle approche qui permet d'apporter une solution au problème de visualisation des données engendré par l'éparsité. En se basant sur les résultats d'une analyse des correspondances multiples (ACM), nous tentons d'atténuer l'effet négatif de l'éparsité en organisant différemment les cellules d'un cube de données. Notre méthode ne cherche pas à réduire l'éparsité mais plutôt à construire un espace de représentation se prêtant mieux à l'analyse et dans lequel les faits du cube sont regroupés. Pour évaluer l'apport de cette nouvelle représentation des données, nous proposons un indice d'homogénéité basé sur le voisinage géométrique des cellules d'un cube. Les différents tests menés nous ont montré l'efficacité de notre méthode.

Mots-clés : ACM, arrangement, cube de données, éparsité d'un cube, espace de représentation, indice d'homogénéité, OLAP, visualisation, voisinage.

1 Introduction

Dans un contexte concurrentiel développé, les entreprises telles que les banques¹ doivent aujourd'hui être capables de prendre des décisions pertinentes, de façon réactive. La mise en place d'un processus décisionnel est alors nécessaire pour gérer une masse de données de plus en plus conséquente. Le stockage et la centralisation de ces données dans un entrepôt constitue un support efficace pour l'analyse de ces dernières. En effet, à partir d'un entrepôt de données, on dispose d'outils permettant de construire des contextes d'analyse multidimensionnels ciblés, appelés communément cubes de données. Ces cubes de données répondent à des besoins d'analyse prédéfinis en amont.

L'analyse en ligne OLAP (On Line Analytical Processing) est un outil basé sur la visualisation permettant la navigation, l'exploration dans ces cubes de données. L'objectif est d'observer des faits, à travers une ou plusieurs mesures, en fonction de différentes dimensions. Il s'agit par exemple d'observer les niveaux de ventes en fonction

¹Nous remercions Michel Rougié, représentant du Crédit Lyonnais, pour les données fournies afin de valider ce travail.

des produits, des périmètres commerciaux (localisations géographiques) et de la période d'achat.

De cette visualisation dépend la qualité d'exploitation des données. Or, différents facteurs peuvent dégrader cette visualisation. D'une part, la représentation multidimensionnelle engendre une éparsité, puisqu'à l'intersection de différentes modalités de dimensions, il n'existe pas forcément de faits correspondants. Cette éparsité peut être accentuée par la considération d'un grand nombre de dimensions (forte dimensionnalité) et/ou d'un grand nombre de modalités dans chacune des dimensions. D'autre part, les modalités des dimensions sont généralement représentées selon un ordre pré-établi (ordre naturel) : ordre chronologique pour les dates, alphabétique pour les libellés. Dans la plupart des cas, cet ordre entraîne une distribution aléatoire des points représentant les faits observés (les cellules pleines) dans l'espace des dimensions.

Dans cet article, nous proposons d'améliorer la visualisation des données dans les cubes. Nous ne diminuons pas l'éparsité du cube comme dans [Niemi *et al.*, 2003], mais à atténuer son effet négatif sur la visualisation, en regroupant les cellules pleines. Pour ce faire, nous proposons d'arranger l'ordre des modalités étant donné que l'ordre initial n'engendre pas forcément une bonne visualisation. Cet arrangement tient compte des corrélations existant entre les faits présents dans l'espace de représentation d'un cube de données. Les corrélations sont fournies par le résultat d'une analyse des correspondances multiples (ACM) appliquée sur les faits du cube.

Ce travail s'inscrit dans une approche générale de couplage entre fouille de données et analyse en ligne. Dans [Messaoud *et al.*, 2005], une réflexion sur l'usage de l'analyse factorielle dans un contexte OLAP a été amorcée. À présent, nous exploitons l'ACM comme un outil d'aide à la construction de cubes de données ayant de meilleures caractéristiques pour la visualisation. En effet, l'ACM construit des axes factoriels qui offrent de meilleurs points de vue du nuage de points des individus.

L'article est organisé comme suit. Dans la section 2, nous repositionnons plus en détail le contexte et les motivations de notre travail. Nous détaillons les différentes étapes de notre approche dans la section 3. Nous présentons dans la section 4 une étude de cas sur un jeu de données bancaires. Dans la section 5, nous donnons un aperçu des travaux connexes au nôtre. Enfin, dans la section 6, nous dressons une conclusion et proposons des perspectives de recherche.

2 Contexte et motivations

Dans un système décisionnel, les données sont organisées selon un modèle, en "étoile" ou en "flocon de neige", dédié à l'analyse et traduisant un contexte d'étude ciblé [Inmon, 1996, Kimball, 1996]. Autour d'une table de faits centrale contenant une ou plusieurs mesures à observer, existent plusieurs tables de dimensions comprenant des descripteurs. Une dimension peut comporter plusieurs hiérarchies impliquant différents niveaux de granularités possibles dans la description de chaque fait. Cette organisation est particulièrement adaptée pour créer des structures multidimensionnelles, appelées "cubes" de données, destinées à l'analyse OLAP. Dans un cube de données, un fait est ainsi identifié par un ensemble de modalités prises par les différentes dimensions. Le fait est observé par une ou plusieurs mesures ayant des propriétés d'additivité plus ou

moins fortes.

La vocation de l'OLAP est de fournir à l'utilisateur un outil visuel pour consulter, explorer et naviguer dans les données d'un cube afin d'y découvrir rapidement et facilement des informations pertinentes. Toutefois, dans le cas de données volumineuses, telles que les données bancaires considérées dans notre étude, l'analyse en ligne n'est pas une tâche facile pour l'utilisateur. En effet, un cube à forte dimensionnalité comportant un grand nombre de modalités, présente souvent une structure éparsée difficile à exploiter visuellement. De plus, l'éparsité, souvent répartie de façon aléatoire dans le cube, altère davantage la qualité de la visualisation et de la navigation dans les données.

Prenons l'exemple de la figure 1 qui présente un cube de données bancaires à deux dimensions : les localités géographiques des agences (L_1, \dots, L_8) et les produits de la banque (P_1, \dots, P_{12}). Les cellules grisées sur la figure sont pleines et représentent la mesure de faits existants (chiffres d'affaires, par exemple) alors que les cellules blanches sont vides et correspondent à des faits inexistant (pas de mesures pour ces croisements de modalités). D'après la figure 1, la répartition des cellules pleines dans la représentation (a) ne se prête pas facilement à l'interprétation. En effet, visuellement, l'information est éparpillée (d'une façon aléatoire) dans l'espace de représentation des données. En revanche, dans la représentation (b), les cellules pleines sont concentrées dans la zone centrale du cube. Cette représentation offre des possibilités de comparaison et d'analyse des valeurs des cellules pleines (les mesures des faits) plus aisées et plus rapides pour l'utilisateur.

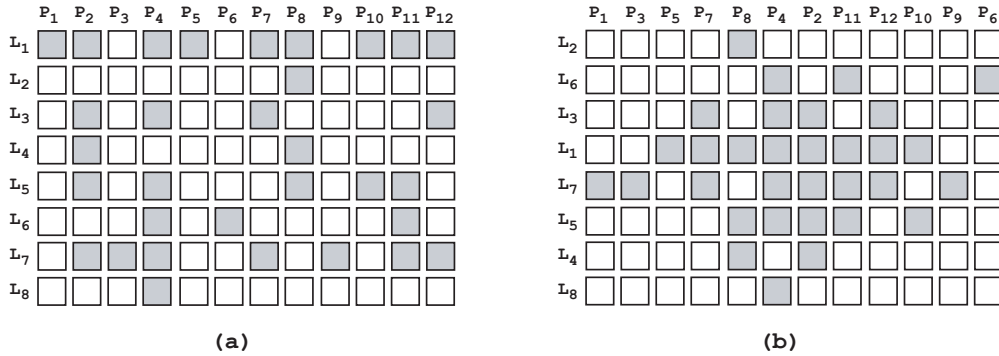


FIG. 1 – Exemple de deux représentations d'un espace de données

Notons que les deux représentations de la figure 1 correspondent au même cube de données. La représentation (b) est obtenue par simples permutations de lignes et de colonnes de la représentation (a). Dans la plupart des serveurs OLAP, les modalités d'une dimension sont présentées selon un ordre arbitraire. En général, cet ordre est alphabétique pour les libellés des modalités et chronologique pour les dimensions temporelles. Malheureusement, dans le cas des cubes éparsés et volumineux, ce choix entraîne des représentations de données inadaptées à l'analyse, voire même difficilement exploitables, comme c'est le cas de la représentation (a) de la figure 1.

La composante visuelle de l'OLAP est primordiale dans un processus décisionnel. En effet, de la qualité et de la clarté de celle-ci dépendent les orientations de l'utilisateur dans son exploration du cube. Ceci détermine la qualité des résultats finaux de l'analyse en ligne. En se basant sur notre idée de l'arrangement des modalités des dimensions illustrée dans l'exemple précédent, nous proposons une méthode permettant à l'utilisateur d'améliorer automatiquement la qualité de la représentation des données. Nous souhaitons produire une meilleure visualisation homogénéisant au mieux le nuage des faits (cellules pleines) et mettant en avant des points de vue intéressants pour l'analyse.

Notre idée d'arrangement consiste à rassembler géométriquement les cellules pleines dans l'espace de représentation des données. Dans ce travail, nous ne cherchons pas à diminuer l'éparsité du cube, mais à l'organiser de manière intelligente pour atténuer l'impact négatif sur la visualisation qu'elle engendre. Nous évaluons l'organisation des données de notre méthode par un indice de qualité de la représentation des données que nous définissons dans la section suivante.

Pour des raisons de complexité de traitements, nous avons exclu la recherche d'un optimum global, voire même local, de l'indice de qualité selon une exploration exhaustive des configurations possibles du cube; c'est à dire, toutes les combinaisons des arrangements possibles des modalités des dimensions du cube. En effet, considérons le cas d'un cube à trois dimensions où chaque dimension comporte seulement 10 modalités. Le nombre de configurations possibles pour ce cube est égal à $A_{10}^{10} \times A_{10}^{10} \times A_{10}^{10} = 10! \times 10! \times 10! \simeq 4,7 \cdot 10^{19}$.

Afin de parvenir à un arrangement convenable des modalités du cube, sans passer par une recherche exhaustive d'un optimum, nous choisissons d'utiliser les résultats d'une analyse en correspondances multiples (ACM) [Benzécri, 1969, Lebart *et al.*, 2000]. L'ACM est alors considérée comme une heuristique appliquée à la volée aux données du cube que l'utilisateur cherche à visualiser. Les individus et les variables de l'ACM correspondent respectivement aux faits et aux dimensions du cube. En construisant des axes factoriels, l'ACM fournit une représentation d'associations entre individus et entre variables dans un espace réduit. Ces axes factoriels permettent d'ajuster au mieux le nuage de points des individus et des variables. Dans le cas de notre approche, afin de mieux représenter les données dans un cube, nous proposons d'exploiter les coordonnées de ses modalités sur les axes factoriels. Ces coordonnées déterminent l'ordre d'arrangement des modalités dans les dimensions. Cependant, l'ACM s'applique sur un tableau disjonctif complet obtenu en remplaçant dans le tableau initial chaque variable qualitative par l'ensemble des variables indicatrices des différentes modalités de cette variable.

Dans la section suivante, nous formalisons les étapes de notre approche. Cette formalisation présente la construction du tableau disjonctif complet à partir du cube de données, l'ACM, l'arrangement des modalités des dimensions et l'indice de qualité de la représentation des données.

3 Formalisation

3.1 Notations

Dans la suite de l'article, nous considérons \mathcal{C} un cube de données à d dimensions, m mesures et n faits ($d, m, n \in \mathbb{N}^*$). Nous adoptons les notations suivantes : $D_1, \dots, D_t, \dots, D_d$ représentent les d dimensions de \mathcal{C} .

Pour la clarté de l'exposé, nous supposons que les dimensions ne comportent pas de hiérarchies. Nous considérons que la dimension D_t ($t \in \{1, \dots, d\}$) est un ensemble de p_t modalités qualitatives. On note a_j^t la $j^{\text{ième}}$ modalité de la dimension D_t . Ainsi, l'ensemble des modalités d'une dimension D_t est $\{a_1^t, \dots, a_j^t, \dots, a_{p_t}^t\}$. Soit $p = \sum_{t=1}^d p_t$ le nombre total de toutes les modalités des d dimensions du cube \mathcal{C} .

Une cellule A dans un cube \mathcal{C} est dite pleine (respectivement, vide) si elle contient une mesure d'un fait existant (respectivement, ne contient pas de faits).

3.2 Aplatissement du cube de données

Pour aplatir le cube \mathcal{C} , nous le représentons sous forme bi-dimensionnelle par un tableau disjonctif complet. Pour chaque dimension D_t ($t \in \{1, \dots, d\}$), nous générons une matrice Z_t à n lignes et p_t colonnes. Z_t est telle que sa $i^{\text{ième}}$ ligne contenant $(p_t - 1)$ fois la valeur 0 et une fois la valeur 1 dans la colonne correspondant à la modalité que prend le fait i ($i \in \{1, \dots, n\}$). Le terme général de la matrice Z_t s'écrit :

$$z_{ij}^t = \begin{cases} 1 & \text{si le fait } i \text{ prend la modalité } a_j^t \text{ de la dimension } D_t \\ 0 & \text{sinon} \end{cases}$$

En juxtaposant les d matrices Z_t , nous construisons la matrice Z à n lignes et p colonnes. $Z = [Z_1, Z_2, \dots, Z_t, \dots, Z_d]$ est un tableau disjonctif complet qui décrit les d positions des n faits du cube \mathcal{C} par un codage binaire.

3.3 Application de l'ACM

À partir du tableau disjonctif complet Z , nous construisons le tableau symétrique $B = Z'Z$ (Z' désigne la transposée de Z) d'ordre (p, p) , qui rassemble les croisements deux à deux de toutes les dimensions du cube \mathcal{C} . B est appelé tableau de contingence de "Burt" associé à Z .

Soit X la matrice diagonale, d'ordre (p, p) , ayant les mêmes éléments diagonaux que B et des zéros ailleurs. Pour trouver les axes factoriels, nous diagonalisons la matrice $S = \frac{1}{d} Z' Z X^{-1}$ dont le terme général est :

$$s_{jj'} = \frac{1}{dz_{.j'}} \sum_{i=1}^n z_{ij} z_{ij'}$$

Après diagonalisation, nous obtenons $(p - d)$ valeurs propres de S notées λ_α ($\alpha \in \{1, \dots, (p - d)\}$). Chaque valeur propre λ_α correspond à un axe factoriel F_α , de vecteur directeur u_α et vérifiant dans \mathbb{R}^p l'équation :

$$S u_\alpha = \lambda_\alpha u_\alpha$$

Les modalités de la dimension D_t sont projetées sur les $(p - d)$ axes factoriels. Soit φ_α^t le vecteur des projections des p_t modalités de D_t sur F_α . Notons que $\varphi_\alpha^{t'} = [\varphi_{\alpha 1}^t, \dots, \varphi_{\alpha j}^t, \dots, \varphi_{\alpha p_t}^t]$.

Désignons par φ_α le vecteur des p projections des modalités de toutes les dimensions sur l'axe factoriel α . Notons que $\varphi_\alpha' = [\varphi_\alpha^1, \dots, \varphi_\alpha^t, \dots, \varphi_\alpha^p]$ et que φ_α vérifie l'équation :

$$\frac{1}{d}X^{-1}Z'Z\varphi_\alpha = \lambda_\alpha\varphi_\alpha$$

La contribution d'une modalité a_j^t dans la construction de l'axe α est évaluée par :

$$Cr_\alpha(a_j^t) = \frac{z_{.j}^t \varphi_{\alpha j}^t{}^2}{nd\lambda_\alpha}$$

Où $z_{.j}^t = \sum_{i=1}^n z_{ij}^t$ correspond au nombre de faits dans le cube \mathcal{C} ayant la modalité a_j^t (poids de la modalité a_j^t dans le cube).

La contribution d'une dimension D_t dans la construction du facteur α est la somme des contributions des modalités de cette dimension, soit :

$$Cr_\alpha(D_t) = \sum_{j=1}^{p_t} Cr_\alpha(a_j^t) = \frac{1}{nd\lambda_\alpha} \sum_{j=1}^{p_t} z_{.j}^t \varphi_{\alpha j}^t{}^2$$

3.4 Arrangement des modalités du cube

Notre idée consiste à associer chaque dimension initiale D_t à un axe factoriel F_α . Pour cela, nous exploitons les contributions relatives des dimensions dans la construction des axes factoriels.

Pour une dimension D_t donnée, nous cherchons, parmi les axes factoriels F_α , celui qui a été le mieux expliqué par les modalités de cette dimension. Nous cherchons à maximiser la valeur de $\lambda_\alpha Cr_\alpha(D_t)$. Il s'agit donc de chercher l'axe F_{α^*} pour lequel la somme des carrés des projections pondérées des modalités de la dimension D_t est maximale. Nous cherchons l'indice α^* vérifiant l'équation suivante :

$$\lambda_{\alpha^*} Cr_{\alpha^*}(D_t) = \max_{\alpha \in \{1, \dots, p-d\}} (\lambda_\alpha Cr_\alpha(D_t))$$

À partir des coordonnées des p_t projections $\varphi_{\alpha^* j}^t$ des modalités a_j^t sur l'axe F_{α^*} , nous appliquons un tri croissant de ces coordonnées. Ce tri fournit un ordre des indices j selon lequel nous arrangeons les modalités a_j^t de la dimension D_t .

L'intérêt de cet arrangement est de converger vers une répartition des modalités de la dimension suivant l'axe factoriel. Cet arrangement a pour effet de concentrer les cases pleines au centre du cube et d'éloigner les cases vides vers les extrémités. Sans diminuer l'éparsité, cette méthode nous permet néanmoins d'améliorer la répartition des données dans le cube. Pour estimer la qualité de cet arrangement, nous proposons un indice pour évaluer l'homogénéité du cube.

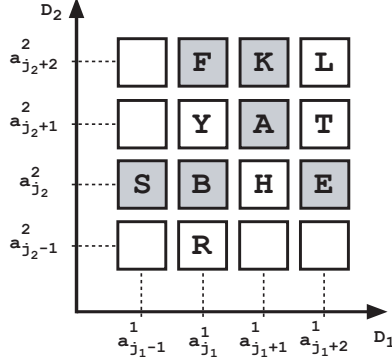


FIG. 2 – Exemple en 2 dimensions de la notion de voisinage des cellules d'un cube de données

3.5 Indice d'homogénéité

Dans cette section, nous proposons un indice permettant de mesurer l'homogénéité de la répartition géométrique des cellules dans un cube. Grâce à cet indice, nous pouvons évaluer le gain induit par l'arrangement des modalités des dimensions. Nous considérons que plus les cellules pleines (ou bien vides) sont concentrées, plus le cube est dit “homogène”.

Une cellule dans un cube représente une ou plusieurs mesures agrégées des faits. Les modalités des dimensions constituent les coordonnées des cellules dans le cube. Soit $A = (a_{j_1}^1, \dots, a_{j_t}^t, \dots, a_{j_d}^d)$ une cellule dans le cube \mathcal{C} , avec $t \in \{1, \dots, d\}$ et $j_t \in \{1, \dots, p_t\}$. j_t est l'indice de la modalité que prend la cellule A pour la dimension D_t .

Nous considérons que toutes les modalités des dimensions D_t sont géométriquement ordonnées dans l'espace de représentation des données selon l'ordre des indices j_t . C'est à dire, la modalité $a_{j_t-1}^t$ précède $a_{j_t}^t$, qui, à son tour, précède $a_{j_t+1}^t$ (voir l'exemple de la figure 2). L'ordre des indices j_t correspond à l'ordre dans lequel sont arrangées dans l'espace les modalités de la dimension D_t . Nous définissons à présent la notion de voisinage pour les cellules d'un cube.

Définition 1 (Cellules voisines) Soit $A = (a_{j_1}^1, \dots, a_{j_t}^t, \dots, a_{j_d}^d)$ une cellule dans un cube \mathcal{C} . La cellule $B = (b_{j_1}^1, \dots, b_{j_t}^t, \dots, b_{j_d}^d)$ est dite voisine de A , notée $B \dashv A$, si $\forall t \in \{1, \dots, d\}$, les coordonnées de B vérifient : $b_{j_t}^t = a_{j_t-1}^t$ ou $b_{j_t}^t = a_{j_t}^t$ ou $b_{j_t}^t = a_{j_t+1}^t$. Exception faite du cas où $\forall t \in \{1, \dots, d\}$ $b_{j_t}^t = a_{j_t}^t$, B n'est pas considérée comme une cellule voisine de A car $B = A$.

Dans l'exemple de la figure 2, la cellule B est voisine de A ($B \dashv A$). Y est aussi voisine de A ($Y \dashv A$). En revanche, les cellules S et R ne sont pas voisines de A . Ceci nous ramène à définir le voisinage d'une cellule.

Définition 2 (Voisinage d'une cellule) Soit A une cellule du cube \mathcal{C} , nous définissons

le voisinage de A , noté $\mathcal{V}(A)$, par l'ensemble de toutes les cellules B de \mathcal{C} qui sont voisines de A .

$$\mathcal{V}(A) = \{B \in \mathcal{C} \text{ tel que } B \dashv A\}$$

Par exemple, dans la figure 2, le voisinage de la cellule A correspond à l'ensemble $\mathcal{V}(A) = \{F, K, L, Y, T, B, H, E\}$.

Définition 3 (Fonction Δ) Nous définissons une fonction Δ de \mathcal{C} dans \mathbb{N} tel que :

$$\forall A \in \mathcal{C}, \Delta(A) = \sum_{B \in \mathcal{V}(A)} \delta(A, B)$$

Avec δ est une fonction définie comme suit :

$$\begin{aligned} \delta : \mathcal{C} \times \mathcal{C} &\longrightarrow \mathbb{N} \\ \delta(A, B) &\longmapsto \begin{cases} 1 & \text{si } A \text{ et } B \text{ sont pleines} \\ 0 & \text{sinon} \end{cases} \end{aligned}$$

$\Delta(A)$ correspond au nombre de cellules pleines et voisines de A .

En supposant que les cellules grises représentent les cellules pleines dans la figure 2, $\Delta(A) = 4$ puisque F , K , B et E sont les seules cellules qui sont à la fois pleines et voisines de A .

Définition 4 (Indice d'homogénéité brut) Nous définissons l'indice d'homogénéité brut d'un cube \mathcal{C} , noté $IHB(\mathcal{C})$, par la somme de tous les couples de ses cellules qui sont à la fois pleines et voisines.

$$IHB(\mathcal{C}) = \sum_{A \in \mathcal{C}} \Delta(A)$$

Par exemple, l'indice d'homogénéité brut du cube de la figure 2 se calcule comme suit :

$$IHB(\mathcal{C}) = \Delta(F) + \Delta(K) + \Delta(A) + \Delta(S) + \Delta(B) + \Delta(E) = 2 + 2 + 4 + 1 + 2 + 1 = 12$$

La meilleure représentation d'un cube de données correspond au cas où ce dernier est complètement non vide. C'est à dire, toutes ses cellules sont pleines. Dans ce cas, l'indice d'homogénéité brut est maximal :

$$IHB_{max}(\mathcal{C}) = \sum_{A \in \mathcal{C}} \sum_{B \in \mathcal{V}(A)} 1$$

Définition 5 (Indice d'homogénéité) Nous définissons, l'indice d'homogénéité d'un cube \mathcal{C} , noté $IH(\mathcal{C})$, par le rapport de l'indice de l'homogénéité brut sur celui de l'homogénéité maximale.

$$IH(\mathcal{C}) = \frac{IHB(\mathcal{C})}{IHB_{max}(\mathcal{C})} = \frac{\sum_{A \in \mathcal{C}} \Delta(A)}{\sum_{A \in \mathcal{C}} \sum_{B \in \mathcal{V}(A)} 1}$$

Dimension	Nombre de modalités	Description
D_1 : catégorie socio-professionnelle	$p_1 = 58$	profil professionnel du client
D_2 : produit	$p_2 = 25$	détention de formule(s) qui sont des offres combinées de produits bancaires
D_3 : unité commerciale	$p_3 = 65$	localisations géographiques de vente
D_4 : segment	$p_4 = 15$	potentiel commercial du client
D_5 : âge	$p_5 = 12$	variable discrétisée selon des tranches d'âge de dix ans ([0-10], [11-20], [21-30], etc.)
D_6 : situation familiale	$p_6 = 6$	exemple : marié, divorcé, etc.
D_7 : type client	$p_7 = 4$	origine du client (par exemple, client membre du personnel du Crédit Lyonnais)
D_8 : marché	$p_8 = 4$	une vente réalisée auprès d'un client est faite sur le marché "particulier des professionnels " si le client est artisan ou exerce une profession libérale, etc., ou sur le marché "particulier" sinon

TAB. 1 – Description des dimensions du cube exemple

Après calcul, l'homogénéité maximale du cube exemple de la figure 2 étant égale à 48, l'indice d'homogénéité de ce dernier est donc $IH(\mathcal{C}) = \frac{12}{48} \simeq 14,28\%$

Pour mesurer l'apport de l'arrangement des modalités sur la représentation du cube de données, nous calculons le gain en homogénéité noté g selon la formule :

$$g = \frac{IH(\mathcal{C}_{arr}) - IH(\mathcal{C}_{ini})}{IH(\mathcal{C}_{ini})}$$

où $IH(\mathcal{C}_{ini})$ est l'indice d'homogénéité de la représentation du cube initial et $IH(\mathcal{C}_{arr})$ est celui de la représentation arrangée selon notre méthode. Notons que quelle que soit la représentation initiale du cube, l'arrangement fourni en sortie par notre méthode est identique puisque l'ACM n'est pas sensible à l'ordre des variables données en entrée.

4 Étude de cas

Pour tester et valider l'approche que nous proposons, nous utilisons un jeu de données bancaires extrait du système d'information du *Crédit Lyonnais*. À partir de ces données, nous avons construit un contexte d'analyse (cube de données). Un fait du cube correspond au comportement d'achat d'un client. Nous disposons dans ce cube de $n = 311\,959$ comportements de clients mesurés par le produit net bancaire (M_1) et le montant des avoirs (M_2). Le tableau 1 détaille la description des dimensions considérées pour observer ces mesures.

Espaces de représentation multidimensionnels dédiés à la visualisation

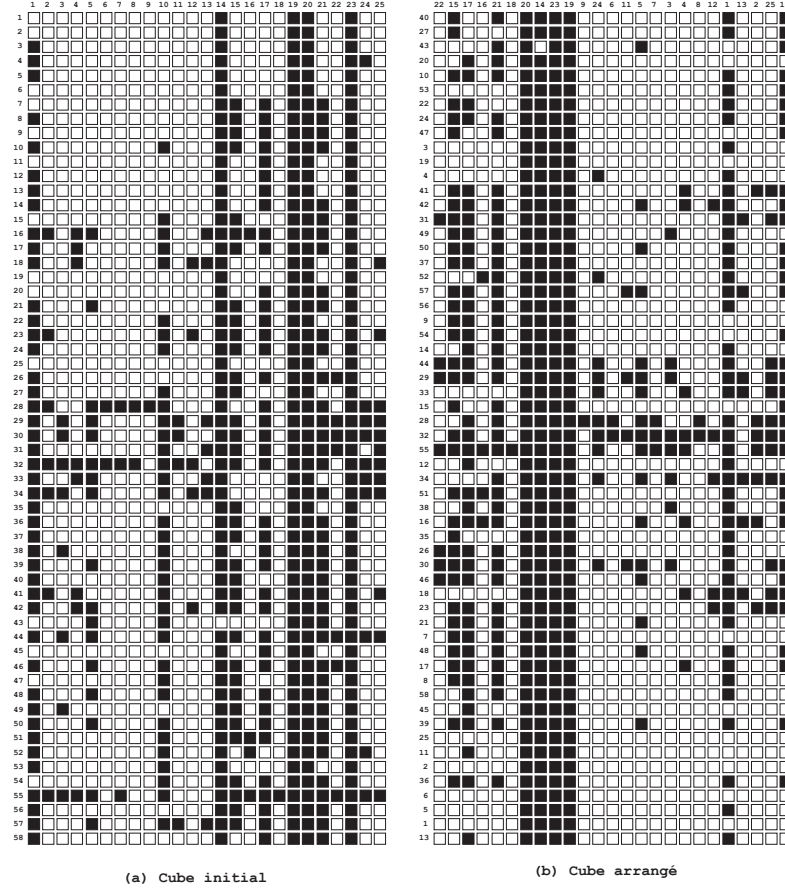


FIG. 3 – Le cube de données avant et après arrangement des modalités

Pour rendre plus claire la suite de notre exposé, notre étude de cas porte sur un cube à deux dimensions ($d = 2$) : la dimension “catégorie socio-professionnelle” (D_1) et la dimension “produit” (D_2). La mesure observée est “le montant des avoirs”. Nous générons les matrices Z_1 et Z_2 selon un codage binaire disjonctif des modalités des deux dimensions. Le tableau disjonctif complet $Z = [Z_1, Z_2]$ a $n = 311\,959$ lignes et $p = p_1 + p_2 = 83$ colonnes.

En appliquant l’ACM sur le tableau Z , on obtient $p - d = 81$ axes factoriels F_α . Chaque axe est caractérisé par sa valeur propre λ_α et les contributions apportées par les dimensions : $Cr_\alpha(D_1)$ et $Cr_\alpha(D_2)$. Nous cherchons, pour chaque dimension, l’axe qui est le mieux contribué par cette dernière. Nous obtenons les résultats suivants :

- Pour la dimension D_1 , $\lambda_{45}Cr_{45}(D_1) = \max_{\alpha \in \{1, \dots, 81\}}(\lambda_\alpha Cr_\alpha(D_1))$, avec $\lambda_{45} = 0.5$ et $Cr_{45}(D_1) = 99.9\%$
- Pour la dimension D_2 , $\lambda_1Cr_1(D_2) = \max_{\alpha \in \{1, \dots, 81\}}(\lambda_\alpha Cr_\alpha(D_2))$, avec $\lambda_1 = 0.83$ et $Cr_1(D_2) = 50\%$.

Ainsi, la dimension D_1 est associée à l'axe F_{45} et D_2 à l'axe F_1 . Les modalités de D_1 (respectivement, D_2) sont arrangées suivant l'ordre croissant de leur projections sur F_{45} (respectivement, F_1). Dans la figure 3, nous présentons le résultat de cet arrangement. La représentation (a) correspond à l'arrangement initial du cube selon l'ordre alphabétique des libellés des modalités. La représentation (b) correspond à l'arrangement obtenu par l'ordre croissant des projections des modalités sur les axes factoriels suscités. Pour des raisons de confidentialité, nous masquons les libellés des modalités de chaque dimension ainsi que les valeurs des mesures. Nous remplaçons les libellés par des codes chiffrés et les mesures existantes par des cases noires. Les cases blanches du cube représentent les creux correspondant à des croisements vides. Sur cet exemple, le taux d'éparsité du cube ² est égal à 64%. La valeur de l'indice d'homogénéité est de 17,75% pour la représentation (a) et de 20,60% pour la représentation (b). Nous obtenons donc un gain en homogénéité de 16,38% par rapport à la représentation initiale du cube.

Nous avons également appliqué notre méthode sur un cube à trois dimensions : "catégorie socio-professionnelle" (D_1), "produit" (D_2) et "âge" (D_5). Ce cube, dont le taux d'éparsité est égal à 87,94%, contient plus de cellules vides comparé au cube précédent. L'arrangement des modalités correspond à l'ordre alphabétique pour D_1 et D_2 , et à l'ordre croissant des tranches d'âge pour D_5 . Le cube initial a un indice d'homogénéité de 5,12%. Le cube arrangé, selon notre méthode, a un indice d'homogénéité de 6,11%. Nous obtenons ainsi un gain de 19,33%.

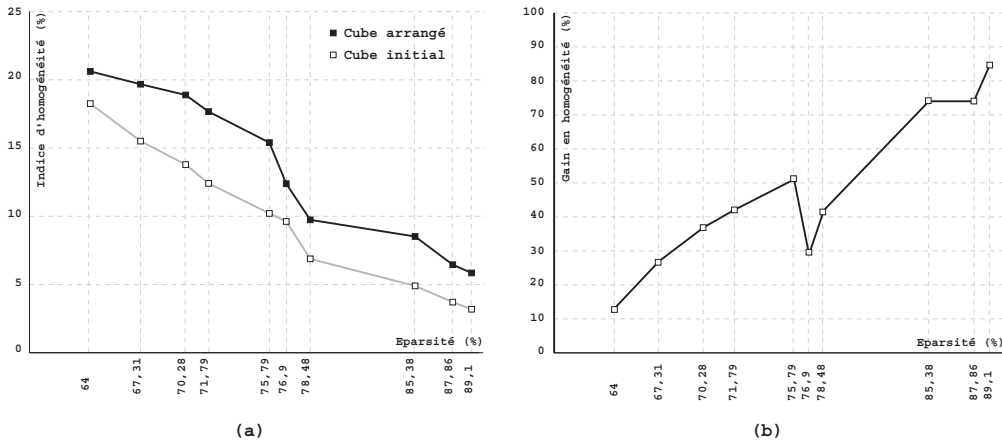


FIG. 4 – Évolutions de l'indice d'homogénéité et du gain en fonction de l'éparsité

Nous avons réalisé une série d'expérimentations de notre méthode sur le premier cube (le cube à deux dimension D_1 et D_2), pour différentes valeurs du taux d'éparsité. Afin de mesurer l'impact de l'éparsité sur notre méthode, nous avons tiré plusieurs échantillons aléatoires à partir de la population du cube initial (les n faits du cube). En variant le taux d'échantillonnage, nous parvenons à faire varier l'éparsité du cube.

²Le taux d'éparsité est égal au rapport entre le nombre de cases vides et le nombre total des cases du cube.

La figure 4 (a) montre l'évolution de l'indice d'homogénéité du cube initial et du cube arrangé en fonction de l'éparsité. Nous remarquons que les valeurs de l'indice sont décroissantes en fonction de l'éparsité du cube. Ceci est naturellement dû à la construction de cet indice qui dépend fortement du nombre de cellules pleines dans le cube. Notons aussi que, quelle que soit l'éparsité, le cube obtenu par arrangement selon notre méthode est toujours de meilleure qualité que le cube initial au sens de notre indice d'homogénéité. Dans tous les cas, nous réalisons un gain en homogénéité lors de l'arrangement du cube.

D'après la figure 4 (b), le gain en homogénéité a une tendance générale croissante en fonction de l'éparsité du cube. En effet, plus le cube est éparse, plus nous avons une meilleure marge de manœuvre pour concentrer les données et les regrouper ensemble autour des axes factoriels de l'ACM.

Notons aussi que le gain en homogénéité, qui est toujours positif, peut fléchir localement (voir figure 4 (b)). Ceci est inhérent à la structure des données. C'est à dire, si les données du cube initial sont déjà dans une représentation homogène, l'application de notre méthode n'apportera pas de gain considérable. En effet, dans ce cas, la méthode n'aura qu'un effet de translation du nuage des fait vers les zones centrales des axes factoriels.

5 Travaux connexes

L'amélioration de l'espace de représentation des données multidimensionnelles dans l'OLAP a fait l'objet de plusieurs travaux de recherche. Rappelons que, dans notre cas, cette amélioration se traduit par la concentration des données autour des axes factoriels d'une ACM. Cela a pour effet de produire une meilleure visualisation homogénéisant au mieux le nuage des faits et mettant en avant des points de vue intéressants pour l'analyse.

Les travaux de recherche qui se sont intéressés à l'étude de l'espace de représentation ont été menés suite à des motivations différentes. Tandis que certains se sont penchés sur des aspects d'optimisation technique (stockage, temps de réponse, etc.), d'autres s'intéressent plutôt à l'aspect de l'analyse en ligne, et particulièrement à la visualisation. Notre travail s'articule davantage autour des seconds travaux. Tout d'abord, nous présentons les travaux ayant traité l'approximation des cubes de données, leur compression et l'optimisation des calculs d'agrégats.

En se basant sur le principe d'approximation par ondelettes (*wavelets*), Vitter *et al.* [Vitter et Wang, 1999] proposent un algorithme pour construire un cube de données compact. L'algorithme proposé fournit des résultats meilleurs que ceux de l'approximation par histogrammes ou par échantillonnage aléatoire [Vitter *et al.*, 1998]. Dans le même ordre d'idées, Barbara et Sullivan [Barbará et Sullivan, 1997] ont proposé l'approche **Quasi-Cube** qui, au lieu de matérialiser la totalité d'un cube, matérialise une partie de ce dernier en se basant sur une description incomplète mais suffisante de ses données. Les données non matérialisées sont ensuite approximées par une régression linéaire.

Une technique de compression basée sur la modélisation statistique de la structure des données d'un cube a été proposée dans [Shanmugasundaram *et al.*, 1999].

Après estimation de la densité de probabilité des données, les auteurs construisent une représentation compacte des données capable de supporter des requêtes d'agrégation. Cette technique n'a de sens que dans le cas de cubes présentant des dimensions continues.

La méthode de compression **Dwarf** proposée dans [Sismanis *et al.*, 2002], réduit l'espace de stockage d'un cube de données. Cette méthode consiste à identifier les n-uplets redondants dans la table de faits. Les redondances de données sont ensuite remplacées par un seul enregistrement. Wang *et al.* [Wang *et al.*, 2002] proposent de factoriser ces redondances par un seul n-uplet de base appelé **BST** (*Base Single Tuple*). À partir du **BST**, les auteurs construisent un cube de données de moindre taille **MinCube** (*Minimal condensed BST Cube*). Cette approche requiert des temps de traitement relativement longs. En vue de remédier à cette limite, Feng *et al.* [Feng *et al.*, 2004a] ont repris l'approche en introduisant une nouvelle structure de données **PrefixCube**. Ils suggèrent de ne plus utiliser tous les **BST** dans la construction du cube mais plutôt de se contenter d'un seul **BST** par dimension. En contre partie, ils proposent l'algorithme **BU-BST** pour la construction d'un cube compressé (*Bottom Up BST algorithm*). Cet algorithme est une version améliorée de l'algorithme **BUC** (*Bottom Up Computation algorithm*) proposé à l'origine dans [Beyer et Ramakrishnan, 1999].

Lakshmanan *et al.* [Lakshmanan *et al.*, 2002] proposent la méthode **Quotient Cube** pour la compression d'un cube de données en résumant son contenu sémantique et en le structurant sous forme de partitions de classes. La meilleure partition n'est pas seulement celle qui permet de réduire la taille du cube mais aussi celle qui permet de conserver une structure de treillis valide donnant la possibilité de naviguer avec les opérations d'agrégation (*Roll-Up*) et de spécification (*Drill-Down*) dans le cube réduit. Malheureusement, la technique des **Quotient Cube** fournit des structures peu compactes. De plus, ces structures ne sont pas adaptées aux mises à jours des données. Dans [Lakshmanan *et al.*, 2003], Lakshmanan *et al.* proposent une nouvelle version améliorée **QC-Tree** (*Quotient Cube Tree*) qui pallie les limites de la technique des **Quotient Cube**. **QC-Tree** permet de rechercher les structures compactes de données dans un cube, d'extraire et de construire les cubes intéressants à partir des données mises à jour.

Feng *et al.* [Feng *et al.*, 2004b] proposent la méthode **Range CUBE** pour la compression des cubes en se basant sur les corrélations entre les cellules du cube. Cette approche consiste à créer un arrangement des cellules d'un cube selon un certain formalisme d'appartenance introduit dans les nœuds du treillis du cube original. Cet arrangement permet de produire une nouvelle structure du cube plus compacte et moins coûteuse en stockage et en temps de réponse.

Ross et Srivastava [Ross et Srivastava, 1997] traitent le problème de l'optimisation du calcul d'agrégats dans les cubes de données éparées. Les auteurs proposent l'algorithme **Partitioned-Cube** qui partitionnent les relations entre les données d'un cube en plusieurs fragments de façon à ce qu'ils tiennent en mémoire centrale. Cette mesure permet de réduire le coût des entrées/sorties. Les fragments de données sont ensuite traités indépendamment, un par un, afin de calculer les agrégats possibles et de générer des sous-cubes de données. Cette notion de fragment est reprise dans les travaux de Li *et al.* [Li *et al.*, 2004]. Leur méthode, appelée **Shell Fragment**, partitionne un ensemble

de données de forte dimensionnalité en sous-ensembles disjoints de données de dimensionnalités moins importantes appelés “*fragments*”. Pour chaque fragment est calculé un cube de données local. Les identifiants des n -uplets participant à la construction de cellules non vides dans un fragment sont enregistrés. Ces identifiants sont utilisés pour lier différents fragments et reconstruire de petits cubes (cuboïdes) nécessaires à l'évaluation d'une requête. Le cube de données de départ est assemblé via ces fragments.

Enfin, citons les travaux de Choong *et al.* [Choong *et al.*, 2004, Choong *et al.*, 2003] qui ont une motivation similaire à la nôtre. Les auteurs utilisent les règles floues (combinaison d'un algorithme de règles d'association et de la théorie des sous-ensembles flous) afin de faciliter la visualisation et la navigation dans l'espace de représentation des cubes de données. Leur approche, consiste à identifier et à construire des blocs de données similaires au sens de la mesure du cube. Cependant, cette approche ne prend pas en compte le problème d'éparsité du cube. De plus, elle se base sur le comptage du nombre d'occurrences des mesures où ces dernières sont considérées comme des nombres entiers.

6 Conclusion et perspectives

Dans cet article, nous avons proposé une nouvelle approche apportant une solution au problème de la visualisation des données dans un cube éparsé. Sans réduire l'éparsité, nous cherchons à organiser l'espace multidimensionnel des données afin de regrouper géométriquement les cellules pleines dans un cube. La recherche d'un arrangement optimal du cube est un problème complexe et coûteux en temps de calcul. Nous avons choisi d'utiliser les résultats de l'ACM comme heuristique pour réduire cette complexité. Notre approche consiste à arranger les modalités des dimensions d'un cube, selon les besoins d'analyse de l'utilisateur, en fonction des résultats fournis par l'ACM. Pour évaluer l'apport de cette nouvelle représentation de données, nous avons proposé un indice d'homogénéité basé sur le voisinage. La comparaison des valeurs de l'indice entre les représentations initiale et arrangée du cube nous permet d'évaluer l'efficacité de notre approche. Les différents tests sur notre jeu de données bancaires nous ont montré, que quelle que soit l'éparsité, notre approche est pertinente. Le gain en homogénéité est croissant en fonction de l'éparsité et son amplitude est également inhérente à la structure des données.

Suite à ce travail, plusieurs perspectives sont à prévoir. Tout d'abord, nous devons étudier la complexité de notre méthode. Cette étude doit prendre en compte aussi bien les propriétés du cube (taille, éparsité, cardinalités, etc.) que l'impact de l'évolution des données (rafraîchissement de l'entrepôt de données).

Ensuite, à ce stade de nos travaux, pour appliquer l'ACM, nous tenons seulement compte de la présence/absence des faits du cube dans la construction des axes factoriels. Nous envisageons alors d'introduire les valeurs des mesures comme pondérations des faits (poids des individus de l'ACM). Ceci permettra de construire des axes factoriels qui traduisent mieux la représentation des faits du cube selon leur ordre de grandeur. Dans ce cas, il serait également intéressant d'introduire la notion de distance entre cellules voisines en fonction des valeurs des mesures qu'elles contiennent.

Dans le même ordre d'idées de la présente méthode, nous souhaitons utiliser les

résultats de l'ACM afin de faire émerger des régions intéressantes à l'analyse à partir d'un cube de données initial. En effet, l'ACM permet de concentrer dans les zones centrales des axes factoriels les individus ayant un comportement normal, et d'éloigner ceux ayant des comportements atypiques vers les zones extrêmes. Nous pouvons déjà exploiter les résultats de l'arrangement des modalités du cube dans le cadre de la distinction de régions correspondant à ces comportements caractéristiques.

Nous voulons aussi comparer la visualisation obtenue par notre approche avec celle proposée dans [Chauchat et Risson, 1998]. Cette dernière représente les résultats d'une analyse factorielle sous forme d'un diagramme de Bertin [Bertin, 1981] qui est plus facile à interpréter. L'objectif de cette méthode est de proposer une visualisation optimisée d'un tableau de contingence. Cependant, elle se limite à des tableaux à deux dimensions sans données manquantes et ne peut pas s'appliquer à des cubes à forte dimensionnalité. Notre approche peut être considérée comme une extension de cette méthode concernant la dimensionnalité du cube et de l'éparcité de ses données.

Par ailleurs, la matérialisation des cubes de données permet le pré-calcul et le stockage des agrégats multidimensionnels de manière à rendre l'analyse OLAP performante. Cela requiert un temps de calcul important et génère un volume de données élevé lorsque le cube matérialisé est à forte dimensionnalité. Au lieu de calculer la totalité du cube, il serait judicieux de calculer et de matérialiser que les parties intéressantes du cube (fragments contenant l'information utile). Comme l'information réside dans les cellules pleines, le cube arrangé obtenu par l'application de l'ACM serait un point de départ pour déterminer ces fragments. Ainsi, comme dans [Barbará et Sullivan, 1997], chaque fragment donnera lieu à un cube local. Les liens entre ces cubes permettront de reconstruire le cube initial.

Enfin, dans ce travail, nous avons délibérément omis de préciser l'origine de ces données. Classiquement, ces données peuvent être issues d'un entrepôt de données. Mais nous envisageons d'appliquer cette approche dans un contexte d'entrepôt virtuel. Nous entendons par entrepôt virtuel la construction de cube à la volée à partir de données fournies par un système de médiation. Un enjeu prometteur de notre méthode est donc de pouvoir soumettre à l'utilisateur, dans le contexte de l'entrepôt virtuel, des représentations visuellement intéressantes des cubes de données. Selon cette démarche, l'utilisateur est de plus en plus impliqué dans le processus décisionnel. D'une part, il est à l'origine des données qu'il veut étudier dans la mesure où il interroge le médiateur. D'autre part, il définit les mesures et les dimensions pour la construction de son contexte d'analyse. Notre méthode se charge alors de lui fournir automatiquement une représentation intéressante en arrangeant les modalités des dimensions qu'il choisit d'observer.

Références

- [Barbará et Sullivan, 1997] Daniel Barbará et Mark Sullivan. Quasi-Cubes : Exploiting Approximations in Multidimensional Databases. *SIGMOD Record*, 26(3) :12–17, 1997.

- [Benzécri, 1969] Jean Paul Benzécri. Statistical analysis as a tool to make patterns emerge from data. In ed.) Academic Press (S. Watanabe, editor, *Methodologies of Pattern Recognition*, pages 35–60, New York, 1969.
- [Bertin, 1981] Jacques Bertin. *Graphics and Graphic Information Processing*. de Gruyter, New York, 1981.
- [Beyer et Ramakrishnan, 1999] Kevin Beyer et Raghu Ramakrishnan. Bottom-Up Computation of Sparse and Iceberg CUBEs. In *Proceedings of ACM SIGMOD Record*, pages 359–370, 1999.
- [Chauchat et Risson, 1998] Jean Hugues Chauchat et Alban Risson. *BERTIN's Graphics and Multidimensional Data Analysis*, pages 37–45. Visualization of Categorical Data. Academic Press., 1998.
- [Choong *et al.*, 2003] Yeow Wei Choong, Dominique Laurent, et Patrick Marcel. Computing Appropriate Representations for Multidimensional Data. *Data & Knowledge Engineering Journal*, 45(2) :181–203, 2003.
- [Choong *et al.*, 2004] Yeow Wei Choong, Anne Laurent, Dominique Laurent, et Pierre Maussion. Résumé de cube de données multidimensionnelles à l'aide de règles floues. In *Revue des Nouvelles Technologies de l'Information*, editor, *4èmes Journées Francophones d'Extraction et de Gestion des Connaissances (EGC 04)*, volume 1, pages 95–106, Clermont-Ferrand, France, Janvier 2004.
- [Feng *et al.*, 2004a] Jianlin Feng, Qiong Fang, et Hulin Ding. PrefixCube : Prefix-sharing Condensed Data Cube. In *Proceedings of the 7th ACM international workshop on Data warehousing and OLAP (DOLAP 04)*, pages 38–47, Washington D.C., U.S.A., November 2004.
- [Feng *et al.*, 2004b] Ying Feng, Divyakant Agrawal, Amr El Abbadi, et Ahmed Metwally. Range CUBE : Efficient Cube Computation by Exploiting Data Correlation. In *Proceedings of the 20th International Conference on Data Engineering*, pages 658–670, 2004.
- [Inmon, 1996] W. H. Inmon. *Building the Data Warehouse*. John Wiley & Sons, 1996.
- [Kimball, 1996] Ralph Kimball. *The Data Warehouse toolkit*. John Wiley & Sons, 1996.
- [Lakashmanan *et al.*, 2002] Laks V.S. Lakashmanan, Jian Pei, et Jiawei Han. Quotient Cube : How to Summarize the Semantics of a Data Cube. In *Proceedings of International Conference of Very Large Data Bases, VLDB'02*, 2002.
- [Lakshmanan *et al.*, 2003] Laks V.S. Lakshmanan, Jian Pei, et Yan Zhao. QC-Trees : An Efficient Summary Structure for Semantic OLAP. In ACM Press, editor, *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, pages 64–75, 2003.
- [Lebart *et al.*, 2000] Ludovic Lebart, Alain Morineau, et Marie Piron. *Statistique exploratoire multidimensionnelle*. Dunod, Paris, 3^e édition edition, 2000.
- [Li *et al.*, 2004] Xiaolei Li, Jiawei Han, et Hector Gonzalez. High-Dimensional OLAP : A Minimal Cubing Approach. In *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB 2004)*, pages 528–539, August 2004.

- [Messaoud *et al.*, 2005] Riadh Ben Messaoud, Sabine Rabaseda, et Omar Boussaid. L'analyse factorielle pour la construction de cubes de données complexes. In *2ème atelier Fouille de Données Complexes dans un processus d'extraction des connaissances, EGC 05, Paris*, pages 53–56, Janvier 2005.
- [Niemi *et al.*, 2003] Tapio Niemi, Jyrki Nummenmaa, et Peter Thanisch. Normalising OLAP cubes for controlling sparsity. *Data & Knowledge Engineering*, 46 :317–343, 2003.
- [Ross et Srivastava, 1997] Kenneth A. Ross et Divesh Srivastava. Fast Computation of Sparse Datacubes. In *Proceedings of the 23rd International Conference of Very Large Data Bases, VLDB'97*, pages 116–125. Morgan Kaufmann, 1997.
- [Shanmugasundaram *et al.*, 1999] Jayavel Shanmugasundaram, Usama M. Fayyad, et Paul S. Bradley. Compressed Data Cubes for OLAP Aggregate Query Approximation on Continuous Dimensions. In *Proceedings of the fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 223–232, August 1999.
- [Sismanis *et al.*, 2002] Yannis Sismanis, Antonios Deligiannakis, Nick Roussopoulos, et Yannis Kotidis. Dwarf : Shrinking the PetaCube. In *Proceedings of the 2002 ACM SIGMOD international conference on Management of data*, pages 464–475. ACM Press, 2002.
- [Vitter *et al.*, 1998] Jeffrey Scott Vitter, Min Wang, et Bala Iyer. Data cube approximation and histograms via wavelets. In *Proceedings of the 7th ACM International Conferences on Information and Knowledge Management (CIKM'98)*, pages 96–104, Washington D.C., U.S.A., November 1998. Association for Computer Machinery.
- [Vitter et Wang, 1999] Jeffrey Scott Vitter et Min Wang. Approximate Computation of Multidimensional Aggregates of Sparse Data Using Wavelets. In *Proceedings of the 1999 ACM SIGMOD international conference on Management of Data*, pages 193–204, Philadelphia, Pennsylvania, U.S.A., June 1999. ACM Press.
- [Wang *et al.*, 2002] Wei Wang, Hongjun Lu, Jianlin Feng, et Jeffrey Xu Yu. Condensed Cube : An Effective Approach to Reducing Data Cube Size. In *Proceedings of the 18th IEEE International Conference on Data Engineering (ICDE'02)*, 2002.

Summary

In decision-support systems, the visual component is important for On Line Analysis Processing (OLAP). In this paper, we propose a new approach that faces the visualization problem due to data sparsity. We use the results of a Multiple Correspondence Analysis (MCA) to reduce the negative effect of sparsity by organizing differently data cube cells. Our approach does not reduce sparsity, however it tries to build relevant representation spaces where facts are efficiently gathered. In order to evaluate our approach, we propose an homogeneity criterion based on geometric neighborhood of cells. The obtained experimental results have shown the efficiency of our method.